

Reproduction and Extension of “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”

Hunner Markus
01503441

Meier Ronja
12433721

Steinegger Benno
12117772

Contents

1	TL;DR	2
2	ModernBERT	2
2.1	Model Changes	2
3	Recreating Benchmarks	2
3.1	GLUE	2
3.2	SWAG	4
3.3	SQuAD 1.1	5
4	Extensions	5
4.1	GLUE: Comparing BERT to ModernBERT	5
4.2	SNLI-Benchmark	6
4.3	Multiple-Choice Class/Head for ModernBERT & GPT-2	7
4.4	Randomly initializing the weights	7
5	Encountered Issue	8
5.1	Reproducing GLUE Experiments	8
6	Distribution of Work	9
7	Appendix	10
7.1	SWAG	10
7.2	SNLI	10

1 TL;DR

All the code can be found in this GitHub repository: <https://github.com/markhun/2024W-DLNLB-BERT>. To get things up and running, just use the Makefile and install the `venv` and install the jupyter kernel and you are ready to go to execute the jupyter notebooks.

2 ModernBERT

ModernBERT [9] is an update to the well-known and established BERT model. This updated model was released just in December 2024 and still relies on the encoder-only architecture against the current trend towards decoder-only generative models. We want to highlight the metaphor from the ModernBERT paper [9]:

We’ve covered this already: encoders are no Ferraris, and ModernBERT is no exception. However, that doesn’t mean it can’t be fast. When you get on the highway, you generally don’t go and trade in your car for a race car, but rather hope that your everyday reliable ride can comfortably hit the speed limit.

2.1 Model Changes

ModernBERT introduces several architectural and training refinements to enhance efficiency and performance. It removes bias terms in all linear layers except the final decoder and disables bias in layer norms. Instead of absolute positional embeddings, it employs Rotary Positional Embeddings (RoPE) for better contextual understanding. Activation functions are updated to GeGLU, a gated-linear variant of BERT’s original GeLU, while normalization is improved with a pre-normalization block and standard layer normalization to stabilize training. The Next Sentence Prediction (NSP) task is removed, and Masked Language Modeling (MLM) is expanded to 30% from the traditional 15%, inspired by MosaicBERT. The tokenizer is a BPE-based model from OLMo, retaining the same special tokens as BERT. Structurally, ModernBERT is deeper but narrower, featuring 22 (12 in original BERT) layers for the base model and 28 (24) for the large version, with a significantly higher context length of 8192 tokens. Learning rate and batch size scheduling have been optimized. Attention mechanisms alternate between global and sliding window attention, incorporating Flash Attention 3 for efficiency. Additionally, padding tokens are removed through unpadding, and training is accelerated using `torch.compile`.

3 Recreating Benchmarks

We did generate a lot of figures and plots while running these benchmarks. The most interesting ones can be found in the appendix 7.

3.1 GLUE

The General Language Understanding Evaluation (GLUE) [8] benchmark is a collection of several natural language understanding tasks, including question answering, sentiment analysis and textual entailment. It consists of nine distinct training datasets, ten validation datasets and eleven test datasets. One of the test datasets (ax) is a diagnostic problem and not part of the benchmark itself. One training dataset (MNLI) provides two validation and test sets, once for entailment and contradiction prediction on in-domain sentence pairs and one for prediction performance

on cross-domain sentence pairs. Table 1 shows an overview over all tasks part of the GLUE benchmark.

Corpus	Train	Val	Task	Metrics	Domain
Single-Sentence Tasks					
CoLA	8.5k	1k	acceptability	Matthews corr.	misc.
SST-2	67k	872	sentiment	acc.	movie reviews
Similarity and Paraphrase Tasks					
MRPC	3.7k	408	paraphrase	acc./F1	news
STS-B	7k	1.5k	sentence similarity	Pearson/Spearman corr.	misc.
QQP	364k	40.4k	paraphrase	acc./F1	social QA questions
Inference Tasks					
MNLI	393k	9.8k/9.8k	NLI	matched acc./mismatched acc.	misc.
QNLI	105k	5.4k	QA/NLI	acc.	Wikipedia
RTE	2.5k	277	NLI	acc.	news, Wikipedia
WNLI	634	71	coreference/NLI	acc.	fiction books

Table 1: Task descriptions and statistics. All tasks are single sentence or sentence pair classification, except STS-B, which is a regression task. MNLI has three classes; all other classification tasks have two. (Table 1 in [8])

In [2] BERT is fine-tuned on all GLUE tasks except for WNLI. The authors decided to not include this task within their evaluation, as it contained adversarial examples at the time of publication. For the tasks MRPC and QQP, only the F1 metric is reported. For the STS-B task, only Spearman correlation is reported.

The labels for all GLUE tasks have never been published. The motivation behind this is to keep the ground truth of the benchmark secret, as to prohibit any publication of results which may have used test data for model training. To evaluate a model on the GLUE benchmark suite, the trained model, together with a pre-print of the paper to be published, has to be sent to the University of New York (NYU) for evaluation. NYU maintains an online leaderboard to allow comparison of all models evaluated on the GLUE benchmark.

Therefore, the main problem in reproducing the GLUE results of [2] is to come up with a labeled test dataset that may be somewhat comparable to original testing data provided by GLUE. To achieve this, we prepare a holdout of the provided training data to be used as test dataset. With one exception: For the MNLI task, GLUE provides two test datasets. One only containing in-domain sentence pairs and another one only containing out-of-domain sentence pairs. To produce two comparable testing datasets, we prepare two holdouts of the respective validation datasets. We did not split a test dataset from the provided validation dataset for other tasks, as most tasks only provide very small validation datasets. Preparing a 10 or 20% holdout from such small validation datasets would result in a test dataset far too small to provide any meaningful results. Fortunately, the validation datasets for the MNLI sets consists of close to 10.000 items and are therefore large enough to extract a new test dataset. No data items included in the newly created test dataset were used during model training or validation. Nevertheless, this means that our reproduced results are strictly speaking not comparable to original results of [2], not only because we evaluate our models on a different testing dataset but also because we had to remove the data items held out for testing from the training or validation dataset.

Table 2 shows the evaluation provided in [2] highlighted in yellow. In the original paper, the result for BERT are compared to three models, which were considered state of the art at the time

System	MNLI (m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Avg.
Current Top 3 Models according to GLUE leaderboard:									
Turing ULR v6	92.5/92.1	76.4	96.7	97.5	73.3	93.1	94.2	93.6	89.9
Vega v1	92.1/91.9	76.7	96.7	97.9	73.8	93.1	94.5	92.4	89.9
Turing NLR v5	92.6/92.4	76.4	97.9	97.6	72.6	93.3	93.8	94.1	90.0
Models compared to in original paper:									
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
Results of original BERT paper:									
BERT _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1
Our reproduction:									
BERT _{BASE}	86.4/84.2	87.7	89.4	94.6	54.8	84.9	85.6	64.7	81.4
BERT _{LARGE}	88.2/85.3	88.2	90.5	94.7	58.8	87.6	88.1	67.5	83.2
Our comparison with ModernBERT:									
ModernBERT _{BASE}	88.4/86.9	88.3	90.6	94.1	58.9	89.0	89.7	76.3	84.7
ModernBERT _{LARGE}	89.2/87.6	88.7	91.5	94.2	60.2	90.6	90.5	83.1	86.2

Table 2: GLUE Test results. The ‘‘Average’’ column is different from the official GLUE score, since the problematic WNLI set was excluded. F1 scores are reported for QQP and MRPC, Spearman correlations are reported for STS-B, Matthew’s correlations are reported for CoLA, and accuracy scores are reported for the other tasks. Original results as provided by [2] are highlighted in yellow.

of publication (2018). To put the results in perspective, we added the results for the current top three models listed on the GLUE leaderboard¹ This shows that many of the tasks included within the GLUE benchmark can no longer be considered particularly hard for current state-of-the-art models, with the top three reaching more than 90 points in the respective metrics for most tasks.

To reproduce the results for BERT_{BASE} and BERT_{LARGE} we fine-tuned the model checkpoints `bert-base-uncased` and `bert-large-uncased` with a learning rate of 2e-5 for five epochs on our modified training data. For evaluation, we choose the best performing epoch checkpoint per GLUE task based on the respective validation dataset performance. As shown in Table 2 this allows us to achieve comparable results to the original paper within all tasks. With the single exception being the QQP task, where our models score more than 18 points above the original results. However, this can anomaly can be attributed to the fact that for the QQP our test dataset is hardly comparable to the original test dataset provided by GLUE (See section 5.1).

3.2 SWAG

The paper ‘‘SWAG: A Large-Scale Adversarial Dataset for Grounded Commonsense Inference’’ [11] introduces SWAG (Situations With Adversarial Generations), a dataset comprising 113,000 multiple-choice questions designed to evaluate a model’s ability to predict subsequent events in various scenarios. Each question presents a context and four possible continuations, with only one being correct. The correct answer may not be directly entailed from the given input sentence but can be derived from the input in combination with commonsense-knowledge. The data stems from video captions. To mitigate annotation artifacts [4, 6] and human biases, from which the model derives the label from the style instead of the content, the authors propose Adversarial

¹See <https://gluebenchmark.com/leaderboard> - Accessed 2025-01-26.

Filtering (AF), a novel method that uses an ensemble of stylistic classifiers to filter out biased data. Empirical results demonstrate that while humans achieve high accuracy (88%) on these questions, models available when the paper was released struggle. Nowadays, models achieve up to 94%².

Unfortunately, the gold labels for the `test`-split are not public, therefore we created our own split by merging the train and validation set and then creating train-validate-test split (60%, 20%, 20% respectively) from the merged dataset. This is one of the reasons why we did not get the same results. We used a batch size of 16 for all experiments, a learning rate of 2e-5 and trained for 3 epochs. In the paper, the authors used an older version of GPT, which is trained on BooksCorpus. We use GPT2, which was introduced after the paper has been released. They also mentioned that they added an additional, trainable vector on top, which is passed through a softmax function for the output. We assume that it is the same as in the `AutoModelForMultipleChoice`-Class.

For fine-tuning BERT on the SWAG dataset, input data is first preprocessed by pairing each context sentence (`sent1`) with all four possible endings (`ending0` – `ending3`), forming multiple-choice options. These pairs are tokenized with truncation and a maximum sequence length of 256, then grouped back into sets of four. The `DataCollatorForMultipleChoice` further processes this data by dynamically padding inputs for uniform batch dimensions, reshaping the flattened inputs back to their original structure, and attaching the corresponding labels as PyTorch tensors. This streamlined pipeline optimizes data handling for multiple-choice classification tasks using Hugging Face’s Trainer.

Model	Validation		Test		Parameters
	Paper	Our	Paper	Our	
BERT-Base	81.6	76.6	-	76.1	110M
BERT-Large	86.6	79.3	86.3	80.8	340M
OpenAI GPT	-	66.8	78.0	66.2	117M
ModernBERT-base	-	77.6	-	77.8	149M
ModernBERT-base _{xavier}	-	77.7	-	78.0	149M
ModernBERT-large	-	82.0	-	82.1	395M
DeBERTa-v3-base	-	83.6	-	83.4	100M

Table 3: All calculated accuracies for the SWAG-Benchmark and the size of the models. DeBERTa-large and not our finetuned DeBERTa-base [5], is the current leader of the SWAG benchmark with over 94%.

3.3 SQuAD 1.1

4 Extensions

4.1 GLUE: Comparing BERT to ModernBERT

As an extension to the evaluation of the original paper we performed the same evaluation on a ModernBERT_{BASE} and ModernBERT_{LARGE} model for each GLUE task. We used a learning

²See <https://leaderboard.allenai.org/swag/submissions/public>

Model	EM		F1	
	Paper	Our	Paper	Our
BERT-Base	80.8	79.9	88.5	87.7
BERT-Large	84.1	83.3	90.9	90.5

Table 4: Calculated EM and F1 score for the SQuAD 1.1 task of each model

rate of $8e-5$ and again choose the best performing epoch checkpoint per GLUE task based on the respective validation dataset performance. As shown in Table 2, ModernBERT_{BASE} usually achieves a performance comparable to BERT_{LARGE}, with the exception of the RTE task where ModernBERT scores roughly 12 to 16 points higher than the original BERT models.

System	MNLI (m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Avg.
BERT _{LARGE}	88.2/85.3	88.2	90.5	94.7	58.8	87.6	88.1	67.5	83.2
ModernBERT _{BASE}	88.4/86.9	88.3	90.6	94.1	58.9	89.0	89.7	76.3	84.7
Speedup (training)									
BERT _{BASE} /ModernBERT _{BASE}	1.38	1.24	1.38	0.99	0.76	1.15	1.18	1.36	1.18

Table 5: GLUE Test results of BERT_{LARGE} and ModernBERT_{BASE} in relation to achieved training time speedup of ModernBERT_{BASE} compared to BERT_{BASE}.

Not only achieves ModernBERT_{BASE} benchmark results comparable to BERT_{LARGE}, but also achieves a better runtime performance compared to BERT_{BASE}. Table 5 shows a comparison of the prediction performance between ModernBERT_{BASE} and BERT_{LARGE} and the achieved speedup of the ModernBERT model compared to BERT_{BASE} in training time. Only for the CoLA task, ModernBERT is slower than the original BERT model. However, this can be explained by the small size of the CoLA task. Training the BERT_{BASE} model on CoLA takes roughly 74 seconds, whereas the ModernBERT variant takes around 97 seconds to train. As ModernBERT consists of more layers and hence more parameters than the original BERT model, any runtime performance advantage of ModernBERT might be overshadowed by a one-time overhead of initialization for such short training times. For the SST-2 task, runtimes are roughly equal. For all other tasks ModernBERT achieves a speedup between 1.15 to 1.38. This means that the smaller ModernBERT model reaches the prediction performance of BERT’s large variant and in addition to way shorter training times than BERT’s base variant.

4.2 SNLI-Benchmark

The SNLI dataset [1] contains a huge corpus (570k) of textual entailment examples, consisting of a text and a hypothesis about the text. The task of the model is to judge whether the hypothesis is entailed, contradicted or neutral to the text. We fine-tuned the pre-trained BERT-base model using the SNLI-dataset. Instead of the Huggingface-Trainer class, we used our own training loop in combination with a custom results class to store loss and accuracy during the different training steps. We achieved a test performance of 90.4% using the BERT-base model. We did also create a custom `NLIClassifier`-class instead of using the `AutoModelForSequenceClassification`-class to have the ability to freeze the pre-trained model, which we ended up not doing.

4.3 Multiple-Choice Class/Head for ModernBERT & GPT-2

As already mentioned in [section 3.2](#), ModernBERT [10] and GPT-2 do not have a MultipleChoice Class/Head. This Python class loads a pre-trained ModernBERT or GPT-2 model, adds dropout for regularization, and includes a classification head to predict the correct answer. The forward function reshapes the input to process all answer choices simultaneously, passes them through ModernBERT/GPT-2, and extracts the hidden state of the first token for ModernBERT and the last token for GPT-2 from each choice. These representations are then passed through the classifier to produce logits for each option, enabling the model to select the most plausible answer. Without this new custom class, the forward function would not be compatible with the structure of the input. In addition, the additional linear classification layer and the softmax function would be missing.

4.4 Randomly initializing the weights

Lastly, we want to explore how weight initialization influences the performance of the model. As Dodge et al. [3] and Sutskever et. al. [7] showcase, weight initialization does influence the performance of (deep) neuronal networks. Our experiment is limited to only one layer with 768 (the size of the last hidden layer of ModernBERT) input weights and one output. We assume that the effect would have been greater if we had more weights to initialize. We used the Xavier uniform initialization schema to initialize the weights. As illustrated in the following plot, it did have a (little) positive effect on the accuracy of the model. The final test accuracy with weight initialization was 77.96 % compared to 77.57 % without weight initialization.

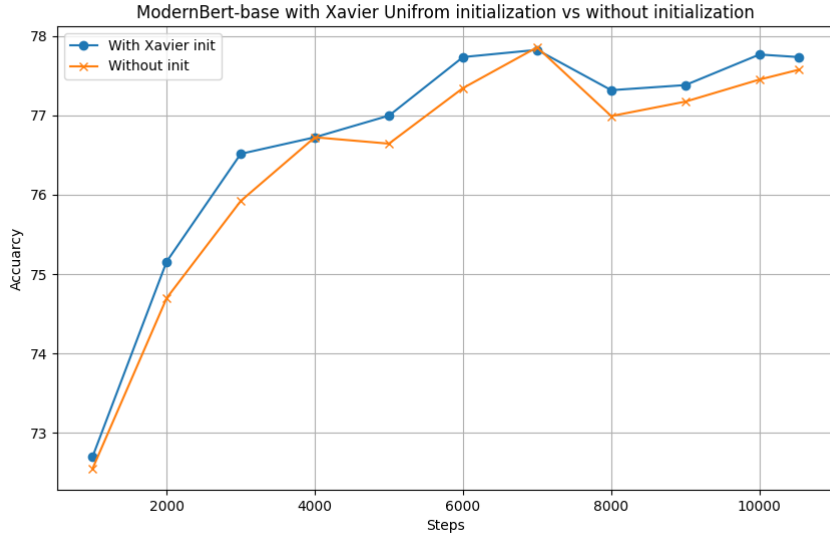


Figure 1: Comparison of validation accuracy between ModernBERT-base with and without weight initialization

Due to computing and time restrictions, we limit this experiment to the SWAG benchmark with the ModernBERT-base model.

5 Encountered Issue

5.1 Reproducing GLUE Experiments

The original paper [2] specifies the exact learning rate used to achieve the reported results on the SWAG and SQuAD benchmarks. For GLUE however the paper just lists several numbers and states the “the best fine-tuning learning rate [was selected]” based on results on the validation data. In turn the authors claimed that they fine-tuned for exactly 3 epochs on each GLUE task, however in the appendix section they state that 2, 3 and 4 epochs are considered *good* for most tasks. We trained for exactly 5 epochs on every GLUE task and then selected the best epoch checkpoint based on performance on the validation dataset. Validation loss plots suggest for some GLUE tasks training for 2 epochs might yield better results than training for 3 or more epochs. As can be seen in Figure 2, validation loss starts to increase very early during the training, suggesting that the model begins to overfit on the training data already after 1 or 2 epochs of fine-tuning.

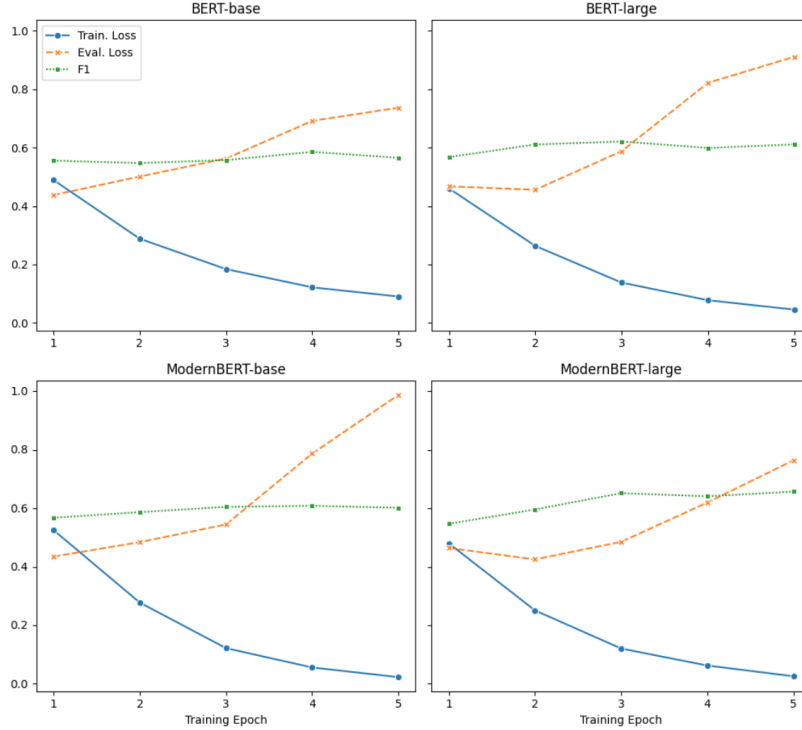


Figure 2: GLUE CoLA fine-tuning - Training loss, Validation loss and F1 score on validation data.

The original table listing the results for GLUE (Table 1 in [2] lists the training set size of every GLUE task within its second row. Unfortunately, the listed numbers do not match the actual dataset sizes provided by GLUE. For example, QNLI is listed with a training set size of “108k” items. However, the current version of GLUE provides a QNLI training dataset with a little less than 105000 items. In the original paper, the authors refrain from training the WNLI task, as it was considered problematic. In the meantime, GLUE updated the datasets for the WNLI task to resolve this issue. The discrepancy in reported training dataset sizes between the original paper

and the current GLUE datasets suggest that also other task datasets were updated. This means that our reproduction of the presented results did not rely on the exact same datasets.

6 Distribution of Work

The following sections were made by Benno Steinegger: SWAG, SNLI-Benchmark, Multiple-Choice Class for ModernBERT & GPT-2, and Randomly initializing the weights. The same split also applies to this report, code and presentation.

The following sections were made by Markus Hunner: 3.1 GLUE, 4.1 GLUE: Comparing BERT to ModernBERT, 5.1 Reproducing GLUE Experiments. The same split also applies to this report, code and presentation.

References

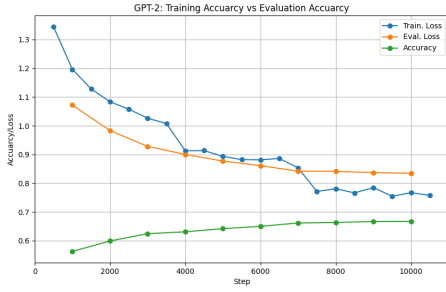
- [1] Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. A large annotated corpus for learning natural language inference, 2015.
- [2] Jacob Devlin. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [3] Jesse Dodge, Gabriel Ilharco, Roy Schwartz, Ali Farhadi, Hannaneh Hajishirzi, and Noah Smith. Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping. *arXiv preprint arXiv:2002.06305*, 2020.
- [4] Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel R. Bowman, and Noah A. Smith. Annotation artifacts in natural language inference data, 2018.
- [5] Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. DeBERTa: Decoding-enhanced BERT with Disentangled Attention, 2021.
- [6] Adam Poliak, Jason Naradowsky, Aparajita Haldar, Rachel Rudinger, and Benjamin Van Durme. Hypothesis only baselines in natural language inference, 2018.
- [7] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In Sanjoy Dasgupta and David McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 1139–1147, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR.
- [8] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding, 2019.
- [9] Benjamin Warner, Antoine Chaffin, Benjamin Clavié, Orion Weller, Oskar Hallström, Said Taghadouini, Alexis Gallagher, Raja Biswas, Faisal Ladhak, Tom Aarsen, et al. Smarter, better, faster, longer: A modern bidirectional encoder for fast, memory efficient, and long context finetuning and inference. *arXiv preprint arXiv:2412.13663*, 2024.

- [10] Benjamin Warner, Antoine Chaffin, Benjamin Clavié, Orion Weller, Oskar Hallström, Said Taghadouini, Alexis Gallagher, Raja Biswas, Faisal Ladhak, Tom Aarsen, et al. Smarter, better, faster, longer: A modern bidirectional encoder for fast, memory efficient, and long context finetuning and inference. *arXiv preprint arXiv:2412.13663*, 2024.
- [11] Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. Swag: A large-scale adversarial dataset for grounded commonsense inference. *arXiv preprint arXiv:1808.05326*, 2018.

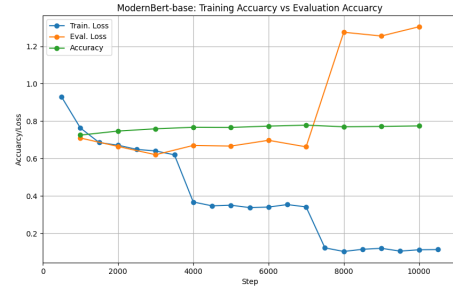
7 Appendix

In this part, we want to include some of the most interesting figures which we generated while fine-tuning. All the figures can be found in the GitHub repo.

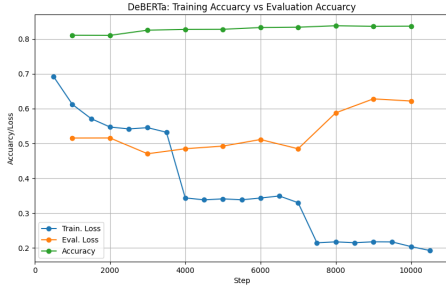
7.1 SWAG



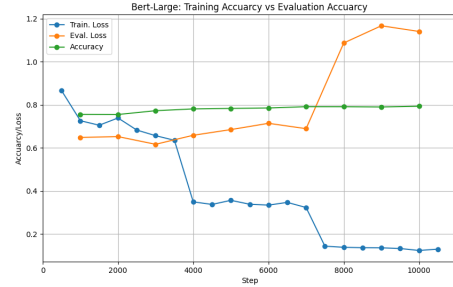
(a) GPT-2 learns without overfitting.



(b) ModernBERT_{BASE} overfits.



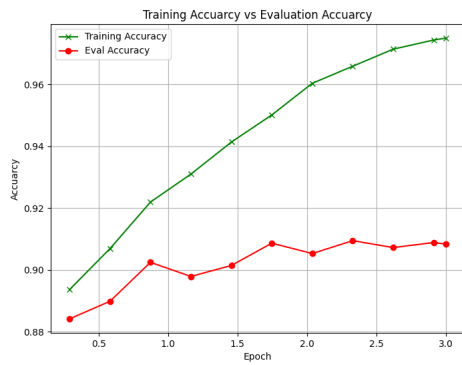
(c) DeBERTa_{BASE} has less loss.



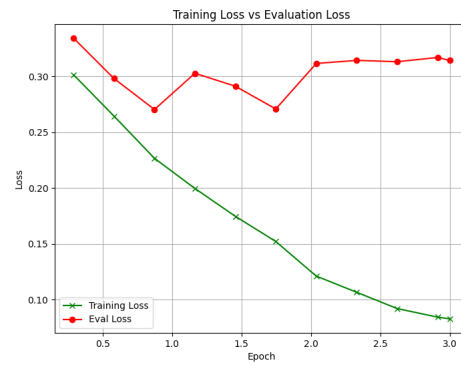
(d) BERT_{LARGE} similar to ModernBERT_{LARGE}.

Figure 3: Performance of different models on SWAG dataset.

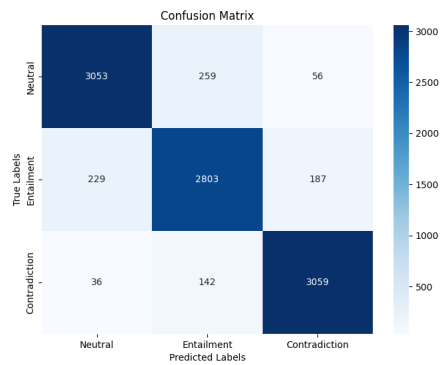
7.2 SNLI



(a) Train vs. val accuracy.



(b) Train vs. val loss.



(c) BERT label classification.

Figure 4: Performance of models on SNLI dataset.